

# On Repairing Structural Problems In Semi-structured Data

Flip Korn, Barna Saha, Divesh Srivastava  
AT&T Labs–Research  
{flip,barna,divesh}@research.att.com

Shanshan Ying  
Nat'l Univ Singapore  
shanshan@comp.nus.edu.sg

## ABSTRACT

Semi-structured data such as XML are popular for data interchange and storage. However, many XML documents have improper nesting where open- and close-tags are unmatched. Since some semi-structured data (e.g., Latex) have a flexible grammar and since many XML documents lack an accompanying DTD or XSD, we focus on computing a syntactic repair via the edit distance.

To solve this problem, we propose a dynamic programming algorithm which takes cubic time. While this algorithm is not scalable, well-formed substrings of the data can be pruned to enable faster computation. Unfortunately, there are still cases where the dynamic program could be very expensive; hence, we give branch-and-bound algorithms based on various combinations of two heuristics, called MinCost and MaxBenefit, that trade off between accuracy and efficiency. Finally, we experimentally demonstrate the performance of these algorithms on real data.

## 1. INTRODUCTION

Semi-structured data provides a flexible representation where data can be nested as a tree and thus is very widely used, from XML documents to JSON data interchange files to annotated linguistic corpora. At the same time, this flexibility makes it prone to errors. A recent study of XML documents on the Web found that 14.6% of them (out of a 180K sample) were not well-formed, the majority of cases due to either open- and close-tag mismatches or missing tags [15]. Such errors are due to multiple factors including manual input [28], dynamically-generated data from faulty scripts [22], mapping and conversion errors (e.g., XML to relational mapping, MS Powerpoint 2007 converted to Powerpoint 2010), and interleaving of multiple sources (e.g., BGPmon pub-sub system which receives XML streams from multiple routers).

Often there is no known grammar associated with the data to test for validity; for example, only 25% of XML documents on the Web have an accompanying DTD or XSD [15]. Inferring one is a notoriously difficult problem [4], often requiring a whole repository rather than a single document, and which for some classes of documents is not even possible [12]. Therefore, most existing work assumes that the document is well-formed and tests validity based

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Articles from this volume were invited to present their results at The 39th International Conference on Very Large Data Bases, August 26th - 30th 2013, Riva del Garda, Trento, Italy.

*Proceedings of the VLDB Endowment*, Vol. 6, No. 9  
Copyright 2013 VLDB Endowment 2150-8097/13/07... \$ 10.00.

on a supplied grammar [24, 23, 5]; exceptions to this are specifically tailored for HTML documents.

In this paper, we consider the problem of repairing an arbitrary semi-structured document into one that is well-formed, based on two variants of well-formedness (discussed below). We believe this problem is in itself interesting for a variety of reasons. First, some existing documents have a very flexible grammar that basically requires only proper nesting. Second, in the absence of a grammar, it may be “safer” to repair based on well-formedness rather than making domain-specific assumptions. Third, since well-formedness is a pre-condition for validity, well-formed repairs may serve as candidates for the user to choose from, similar to the way word processors suggest autocorrection.

While *verifying* well-formedness in semi-structured data can be done straightforwardly, using a stack, in time linear in the size of the document, it is a much more challenging problem to *repair* a malformed document. Some existing tools, such as modern Web browsers, use simple rule-based heuristics to rectify mismatching tags. Perhaps the most common rule, employed by some web browsers such as Internet Explorer, is to substitute a matching close-tag whenever the current close-tag does not match the open-tag on the stack. However, a single extra or missing close-tag is enough to set off a cascade, requiring many close-tags to be replaced (or deleted). Another commonly used rule is to insert a matching close-tag whenever the current close-tag does not match, but this can trigger a similar cascade.

**Example:** Figure 1(a) shows an example XML document of a bibliographic entry that is not well-formed: the `<authors>` open tag does not have a matching close tag; `<affiliation>` occurs out of place and is missing a matching tag; the `</title>` close tag is out of order, occurring after `<authors>`; etc. Figure 1(b) shows the document after the substitution rule-based heuristic is applied, requiring 3 substitutions and 2 insertions. □

We focus on the following types of errors which we believe occur most frequently in practice:

1. Tags may be missing, as it is common to forget to close open tags, and unmatched close tags may occur when new content is added and it is assumed a previous open tag existed.
2. Extraneous tags may be present, perhaps due to not fully deleting tags associated with deleted content.
3. Open and close tags, due to being similar, are sometimes mistaken for each other; and tags of different types may appear in the wrong order or be improperly assigned.

We use standard string edit distance with insertion, deletion and substitution operators as a model for repair [17]. We believe that more complex distance functions including other operations, such

```

<article>
  <title>
    A Relational Model for Large Shared Data Banks
  </title>
  <authors>
    <author>
      <name> E. F. Codd </name>
      IBM <affiliation>
    </author>
  </authors>
</article>

```

(a) original document

```

<article>
  <title>
    A Relational Model for Large Shared Data Banks
  </title>
  <authors>
    <author>
      <name> E. F. Codd </name>
      IBM <affiliation> </affiliation>
    </author>
  </authors>
</article>

```

(b) rule-based repair

Figure 1: An Example and Rule-based Repair

as block moves and swaps, as well as non-uniform weighting, can be folded into our methods but we leave this to future work. Edit distance is used for modeling and correcting errors in many applications from information retrieval to computational biology [30, 20, 27]. In our illustrative example, a well-formed repair with fewest edits is given in Figure 2(a), which has edit distance 2: delete `<authors>` and delete `<affiliation>`.

In our second variant of well-formedness, we take into account that the text embedded within semi-structured documents often follows certain patterns. For example, many XML documents only allow text to occur surrounded by matching open-close tags and require the existence of text between every adjacent matching pair. Thus we consider how to exploit embedded text to aid in finding a more judicious repair via a constrained edit distance function. In our illustrative example, a well-formed repair based on tags and text with edit distance 3 is given in Figure 2(b): delete `<authors>`, insert `<affiliation>` before IBM, and substitute `<affiliation>` after IBM by `</affiliation>`. Note that this repair consists of more edits than for tags only.

Note that it is not always possible to exactly repair to the originally intended well-formed string. In the absence of a grammar, there is inherent ambiguity in what the creator intended. For example, consider the string `<name> E. F. Codd </author>`. Should this be repaired to `<name> E. F. Codd </name>` or `<author> E. F. Codd </author>`? Or even to `<name> E. F. </name> <author> Codd </author>`. It is impossible to know what the original intent was. Furthermore, such ambiguities compound in larger strings, resulting in an explosive number of reasonable possibilities. Since the user may have a (often ill-defined) grammar in mind, our methods can provide multiple repairs in the hope that at least one of these will suffice. But presenting the user with all repairs based on the many ways to resolve these ambiguities can be overwhelming. Instead, we note that the differences between some repairs are syntactically trivial, so we try to consolidate these into representative repairs. For example among the two alternatives `<name> E. F. Codd </name>` and `<author> E. F. Codd </author>` to repair `<name> E. F. Codd </author>`, we canonically choose the former (see more details in Section 3.2). Consolidating multiple repairs by such representatives helps to provide more variety in a small set of repairs returned to the user. For the Tags With Text case, the surrounding text can be exploited to resolve more of these ambiguities. For example, if from a well-formed string such as `<name> E. F. Codd </name>` a tag gets deleted resulting in `<name> E. F. Codd` or `E. F. Codd </name>`, indeed our algorithm will repair it by inserting the deleted tag. Therefore, with a stronger grammar, there are more cues to recover the original string.

There has been much literature on approximate matching of trees which has been applied to finding semantically relevant XML doc-

uments [13, 9, 8, 29, 21, 31]. Unfortunately, none of this work applies to our setting since the input is not well-formed and, therefore, cannot be represented as a tree. However, a good repair should result in a short tree edit distance between the repaired string and the intended error-free string. We use this to show the efficacy of our algorithms in “undoing” errors introduced to a well-formed string in Section 5.2. Recall that for the reasons of ambiguity mentioned above, it is not enough to simply check whether or not the repaired string is the *same as* the intended error-free string. In addition, our experimental evaluation on real XML data with real errors shows that the number of string edit operations is much smaller when using our approach compared to the rule-based heuristics (see Section 5.2). *This effectively establishes the goodness of edit distance for repair.*

Our contributions are as follows. We give a dynamic programming algorithm which computes the (optimal) edit distance. Since this algorithm is cubic in the size of the input, regardless of the number of errors, it does not scale to large documents. Therefore, we present linear-time greedy algorithms that find close to optimal (within 10%) repairs in practice; these methods are significantly more accurate (more than an order of magnitude) than the rule-based heuristics. We also propose branch-and-bound algorithms for when multiple repairs are desired (such as for an autocorrection menu), since the dynamic program and greedy algorithms are geared towards finding a single repair. We present a variety of methods, with various trade-offs in accuracy and running time, whose performance depends on the number of edits rather than the length of the input. We perform a thorough experimental study to investigate these strategies on real data.

## 2. RELATED WORK

While we are not aware of prior work that specifically addresses the problem of repairing malformed data to make it syntactically well-formed, there is some work on repairing XML documents to make them valid with respect to a given DTD [5, 25, 26]. However, these papers all assume the input is already well-formed. It is not clear how the techniques used in these papers, such as computing the tree or graph edit distance between a document and a DTD, can be applied to the problem here where documents are malformed. Some existing tools such as *Beautiful Soup*, *HTML Tidy* and *NekoHTML* allow for malformed HTML input and exploit domain knowledge to make them valid; however, they employ simple rules that can result in verbose repairs.

In the absence of a grammar, a natural approach would be to infer one for making repairs. For example, a few papers have been written about using a repository of XML documents to generate a DTD [11, 3] or XSD [4]. Unfortunately, grammar inference is a notoriously difficult problem, requiring a large number of examples and which is impossible for some classes of grammars when only positive examples are given [12].

```

<article>
  <title>
    A Relational Model for Large Shared Data Banks
  </title>
  <author>
    <name> E. F. Codd </name>
    IBM
  </author>
</article>

```

(a) considering tags only

```

<article>
  <title>
    A Relational Model for Large Shared Data Banks
  </title>
  <author>
    <name> E. F. Codd </name>
    <affiliation> IBM </affiliation>
  </author>
</article>

```

(b) considering tags and text

Figure 2: Two Possible Repairs

The problem of computing the edit distance from a string to a supplied context-free grammar has been studied; since the grammars for our notions of well-formedness can be expressed using a CFG, these existing solutions can be applied. Aho and Peterson [1] gave an  $O(|G|^2 n^3)$  algorithm which was later improved to  $O(|G|n^3)$  by Myers [19], where  $n$  is the length of the input and  $|G| = \sum_{A \rightarrow \alpha \in G} (|\alpha| + 1)$  is the size of the grammar. Our dynamic program runs in  $O(n^3)$  time, independent of grammar size.

It has been shown that a *non-deterministic version* of the language of well-formed bracketed strings is, in terms of parsing, the hardest CFG [14]. It is also known that parsing an arbitrary CFG is at least as hard as boolean matrix multiplication [16]. Therefore, computing the edit distance to a well-formed string in much less than cubic time would be a significant accomplishment.

Verifying well-formedness is a related but much easier problem: it is straightforward to do this using a stack in linear time. The problem is non-trivial, however, on streaming data where trading off accuracy (where distance to well-formedness is measured by Hamming distance) can allow this in sublinear space [18]. Other papers study the problem of validity checking: using a DTD or XML Schema, report if a given input document conforms to the given grammar. Some of these papers (e.g., [24]) perform *strong validation*, checking for well-formedness along with validity, while others (e.g., [23]) perform *weak validation*, assuming the input is already well-formed. In fact, our techniques can be used for preprocessing malformed input to enable application of the latter work.

Finally, we mention that our work fits into the context of data cleaning to satisfy database integrity constraints, including consistency under functional dependencies [2], inclusion dependencies [6] and record matching [10]. These papers all use edit distance as a notion of a minimal cost repair.

### 3. TAGS-ONLY CASE

We assume that the input data is tokenizable by a lexical analyzer and has been preprocessed into a sequence of *brackets*. For example, these brackets could correspond to the open-tags ( $\langle \dots \rangle$ ) and close-tags ( $\langle / \dots \rangle$ ) of an XML document; to the curly braces or square brackets (and accompanying object name) of a JSON file; to a Latex file containing  $\backslash \text{begin}\{\dots\}$  and  $\backslash \text{end}\{\dots\}$ ; etc.<sup>1</sup>

**Definition 1.** The *congruent* of a bracket  $x$  is defined as its symmetric opposite bracket, denoted  $\bar{x}$ . The congruent of a set of brackets  $X$ , denoted  $\bar{X}$ , is defined as  $\{\bar{x} \mid x \in X\}$ .

We assume a bracket namespace is not given *a priori*. Let  $R$  and  $S$  denote the sets of brackets obtained from the two directions (i.e., open and close) after tokenization, respectively. We shall use  $T = R \cup \bar{S}$  to denote the set of *open brackets* and  $\bar{T} = \bar{R} \cup S$  the

<sup>1</sup>We shall ignore all other components in the XML document besides tags for now such as attributes and text and treat empty-element tags as two tags, e.g.,  $\langle a / \rangle$  becomes  $\langle a \rangle \langle / a \rangle$ .

*close brackets*. (Each  $x \in T$  has exactly one congruent  $\bar{x} \in \bar{T}$  and vice versa.)

**Definition 2.** A *match* between two brackets  $x$  and  $y$ , denoted  $x \asymp y$ , occurs when  $x \in T$ ,  $y \in \bar{T}$  and  $y = \bar{x}$ .

Consider a string  $s = s_1 \dots s_n$ , of length  $n$ , over some bracket alphabet  $T$ , that is,  $s \in (T \cup \bar{T})^*$ .

**Definition 3.** A *well-formed string* over some bracket alphabet  $T$  obeys the context-free grammar  $G_T$  with productions  $S \rightarrow SS$ ,  $S \rightarrow \varepsilon$  and  $S \rightarrow aS\bar{a}$  for all  $a \in T$ .<sup>2</sup> So  $|G_T| = 4|T| + 3$ .

**Example:** Let  $T = \{a, b, c\}$ . Then  $ab\bar{b}c\bar{c}\bar{a}$  is a well-formed string, since it can be parsed as  $S \rightarrow aS\bar{a} \rightarrow a(SS)\bar{a} \rightarrow a(bS\bar{b})(cS\bar{c})\bar{a} \rightarrow ab(\varepsilon)\bar{b}c(\varepsilon)\bar{c}\bar{a} \rightarrow ab\bar{b}c\bar{c}\bar{a}$ . However,  $ab\bar{a}\bar{b}$  is not well-formed.  $\square$

**Definition 4.** A *well-formed bracketed language*  $L(G_T)$  over some bracket alphabet  $T$  is the set of strings from  $T \cup \bar{T}^*$  accepted by the grammar  $G_T$  defined above.

**Definition 5.** The *edit distance*  $E(s, s')$  between two strings  $s$  and  $s'$  is the minimum number of insertions, deletions and substitutions needed to transform  $s$  into  $s'$ , where an insertion of  $a$  after position  $i$  transforms  $s_1 \dots s_i s_{i+1} \dots s_n$  to  $s_1 \dots s_i a s_{i+1} \dots s_n$ ; a deletion at position  $i$  transforms  $s_1 \dots s_{i-1} s_i s_{i+1} \dots s_n$  to  $s_1 \dots s_{i-1} s_{i+1} \dots s_n$ ; and a substitution to  $a$  at position  $i$  transforms  $s_1 \dots s_{i-1} s_i s_{i+1} \dots s_n$  to  $s_1 \dots s_{i-1} a s_{i+1} \dots s_n$ .

**Definition 6.** The *Bracketed Language Edit Distance Problem*, given string  $s$ , is to find  $\arg \min_{s'} E(s, s')$  such that  $s' \in L(G_T)$ .

We shall henceforth use the term *edit distance* of a string to mean the edit distance from the string to a well-formed repair.

**Example:** The edit distance of  $ab\bar{a}\bar{b}$  to a well-formed string is 2. For example, it can be changed to  $a\bar{a}$  using 2 deletions and  $ab\bar{b}\bar{a}$  using 2 substitutions.  $\square$

### 3.1 Optimal Solution

Algorithm 1 presents pseudocode for a dynamic program to compute the edit distance  $\min_{s'} E(s, s')$  to a well-formed string  $s'$ . The algorithm runs in  $O(n^3)$  time and requires  $O(n^2)$  space, where  $n$  is the string length and is based on applying the different grammar productions to the input. Given some substring  $s_i \dots s_j$ ,  $j > i$ , either  $s_i$  and  $s_j$  could be edited to matching brackets (the  $S \rightarrow aS\bar{a}$  production) or else the string could be broken into two adjacent well-formed substrings (the  $S \rightarrow SS$  production). Let  $B[i, j]$  be the cost of editing  $s_i$  and  $s_j$  to match. When  $j = i$ ,

<sup>2</sup>Some instances of well-formedness additionally require that the document is nested within a single open-close pair, i.e.,  $s \in T(T \cup \bar{T})^* \bar{T}$ , but we dispense with this for simplicity.

$B[i, j] = B[i, i] = 1$ , also  $C[i, j] = C[i, i] = 1$ . Hence, the recurrence is

$$C[i, j] \leftarrow \min(B[i, j] + C[i+1, j-1], \min_{i \leq k \leq j-1} C[i, k] + C[k+1, j])$$

The following claim establishes the correctness of the recursion and hence, also of Algorithm 1.<sup>3</sup>

**CLAIM 1.** *Algorithm 1 correctly finds the edit distance, given string  $s$ , such that it is accepted by  $G_T$ .*

*Proof Sketch.* Consider any substring of length  $l$ ,  $s_i$ . The minimum edit distance is 1, which can be achieved either by deleting  $s_i$  (by rule  $S \rightarrow \varepsilon$ ), or by inserting matching open/close bracket. For substrings of length 2,  $s_i s_{i+1}$ ,  $B[i, i+1]$  is the computed edit distance and  $C[i, i+1] = B[i, i+1]$ . These serve as the base cases. Suppose, by induction hypothesis, Algorithm 1 correctly computes minimum edit distance for all substrings of length at most  $l$ . We now take any substring of length  $l+1$ . W.l.o.g, let it be  $s_1 s_2 \dots s_{l+1}$ . Consider  $s_1$ . An optimal algorithm has the following options.

- (1) Deletes  $s_1$  and the minimum edit distance is  $1 + C[2, l+1]$ ,
- (2) Matches  $s_1$  to some  $s_j$ ,  $j > 1$ , possibly by substitution. There are again few subcases here.

(a)  $j = l+1$  and  $s_1 \asymp s_{l+1}$ . In that case, the minimum edit distance is  $C[1, l+1] = C[2, l]$ .

(b)  $j = l+1$ ,  $s_1 \in T$  and substituted to match  $s_{l+1}$ . In that case, the minimum edit distance is  $C[1, l+1] = 1 + C[2, l]$ .

(c)  $j = l+1$ ,  $s_{l+1} \in \bar{T}$  and substituted to match  $s_1$ . In that case, the minimum edit distance is  $C[1, l+1] = 1 + C[2, l]$ .

(d)  $j = l+1$ ,  $s_1 \in \bar{T}$  and  $s_{l+1} \in T$ , both deleted or substituted to match each other. In that case, the minimum edit distance is  $C[1, l+1] = 2 + C[2, l]$ .

(e)  $j < l+1$ . In that case, the minimum edit distance is  $C[1, l+1] = C[1, j] + C[j+1, l+1]$ .

The above options are exhaustive, Algorithm 1 considers all the options, computes edit distance correctly in all these cases and returns the minimum. Hence, Algorithm 1 computes minimum edit distance for any substring of length  $l+1$  correctly, therefore, by induction, the proof is established.  $\square$

For ease of exposition, we do not show how to construct an  $s'$ . A single minimum cost repair can be constructed from the dynamic programming tableau straightforwardly. However, constructing multiple repairs having minimum cost is non-trivial. We defer this discussion until Section 3.2.

---

### Algorithm 1 Dynamic Program for Edit Distance

---

**Require:** tokenized string  $s = s_1 \dots s_n$   
**Ensure:** edit distance  $\min_{s'} E(s, s')$  where  $s'$  is well-formed

- 1: **for all**  $\ell$  from 1 to  $n-1$  **do**
- 2:   **for all**  $i$  from 1 to  $n-\ell$  **do**
- 3:      $j \leftarrow i + \ell$
- 4:      $C[i, j] \leftarrow B[i, j] + C[i+1, j-1]$
- 5:     **for all**  $k$  from  $i$  to  $j-1$  **do**
- 6:        $C[i, j] \leftarrow \min(C[i, j], C[i, k] + C[k+1, j])$
- 7:     **end for**
- 8:   **end for**
- 9: **end for**
- 10: **return**  $C[1, n]$

---

The cubic growth in running time as a function of string size becomes a problem for large strings. Luckily, one can often prune

<sup>3</sup>Full version with proofs at <http://www.cs.umd.edu/~barna/vldb13-xmlrepair.pdf>

away well-formed substrings to speed up the running time. Using a stack, this can be done straightforwardly in linear time by recursively finding matching adjacent pairs and removing them. We note, however, that this may eliminate some candidate repairs from consideration. For example, given the string  $a\bar{a}\bar{b}\bar{b}a$ , the only optimal repair with pruning that can be found is  $a\bar{a}a\bar{a}$  whereas  $abb\bar{a}$  also has an edit distance of 1. Nonetheless, there exists at least one repair of the pruned string with edit distance equal to that optimal for the original string.

**CLAIM 2.** *Well-formed substrings removal preserves the edit distance. That is, the edit distance of the pruned string is the same as the edit distance for the full string.*

*Proof Sketch.* We do an induction on the number of edits. For strings with edit distance 0, well-formed substring pruning returns the empty string with edit distance 0. This serves as the basis. Considering the claim to be true for all strings with minimum edit distance less than  $d$ , we next take any string  $s$  with minimal edit distance  $d$ . Consider an optimal algorithm that defers doing the first edit as much as possible without affecting the optimality, and let the  $d$  edit positions be  $p[1] < p[2] < \dots < p[d]$ . Also, let  $t$  be the string that results from  $s$  after removing the well-formed substrings. Let the prefix in  $t$  corresponding to  $s_1 s_2 \dots s_{p[1]}$  be  $t_1 t_2 \dots t_q$ . Clearly,  $t_1, \dots, t_{q-1}$  correspond only to open brackets.

Case 1: If  $s_{p[1]}$  is not part of any well-formed substring then  $t_q$  corresponds to  $s_{p[1]}$ . Consider performing the same edit operation as the optimal algorithm at  $s_{p[1]}$  and also at  $t_q$ . The resultant string  $s'$  after the edit at  $s_{p[1]}$  in  $s$ , has edit distance  $d-1$ . If well-formed substrings are removed from  $s'$ , and also from  $t$  after the edit at  $t_q$  to get  $t'$ , then they both return the same processed string. By the induction hypothesis,  $s'$  and  $t'$  have the same edit distance. Edit distance of  $s$  and  $t$  are one more than edit distance of  $s'$  and  $t'$  respectively, hence, both  $s$  and  $t$  have the same edit distance.

Case 2: Otherwise  $s_{p[1]}$  is part of a well-formed substring  $w$ . There are few cases to consider based on whether  $s_{p[1]} \in T$  or  $s_{p[1]} \in \bar{T}$ . The main idea is to show that in all these cases, there exists an alternate edit script such that  $s_{p[1]}$  is not part of any well-formed substrings. Then by the same argument as in Case 1, the proof follows.  $\square$

There are instances where well-formed substring pruning will not be very effective; for example, consider the string  $abcded\bar{c}\bar{b}a$ . Here edit distance is 1, but since there is no well-formed substring, nothing can be eliminated. We investigate its effectiveness on real data in Section 5.

## 3.2 Incremental Approach

The dynamic program presented in the previous section has two deficiencies. The first is that it has the same running time regardless of how many errors exist in the input string; that is, its best-case running time is as slow as the worst-case. Second, it can extract a single edit script associated with the edit distance found but does not provide a natural way of enumerating multiple repairs. Here we describe branch-and-bound strategies, with various trade-offs between accuracy and running time, that are affected only by the number of errors, not the length of the string, and are capable of incrementally reporting repairs. Our algorithms are based on various combinations of greedy heuristics. All of our methods maintain a stack that, at any point, contains open brackets remaining to be matched with close brackets from the string.

As a warm-up, we consider the case where  $|T| = 1$ . Here it turns out we can apply recursive matching of adjacent pairs to obtain a sequence of zero or more elements from  $\bar{T}$  followed by a sequence of zero or more elements from  $T$ , that is,  $\bar{a}^* a^*$ . Then

the minimum cost repair, for the close bracket and open bracket substrings separately, is obtained by applying substitutions to make adjacent pairs match and, if a singleton remains, delete it. So if the pruned string is  $\bar{a}^i a^j$ , the resulting edit distance is  $\lceil i/2 \rceil + \lceil j/2 \rceil$ . Hence, for  $|T| = 1$  the edit distance can be computed in  $\Theta(n)$  time. Clearly the same edit distance can be obtained via many different matchings, precisely  $p(i/2) \times p(j/2)$  when  $i$  and  $j$  are even, where  $p(k)$  is the partition function denoting the number of ways to write  $k$  as a sum of positive integers. Rather than enumerating all these possibilities, a single canonical form such as only the minimum- or maximum-depth nesting shall be reported. For example, a malformed string such as  $aaaaaa$  will be repaired to  $aaa\bar{a}\bar{a}$  or  $\bar{a}\bar{a}\bar{a}\bar{a}$  but not to  $a\bar{a}\bar{a}\bar{a}$ .

When  $|T| \geq 2$ , repairs of these two scenarios are similar to the  $|T| = 1$  case: adjacent pairs are examined and the appropriate substitutions are made to make these adjacent pairs match. There is an additional type of error than can occur besides these two: empty stack with remaining close brackets and a non-empty stack when the string has terminated. There could be an open bracket of one type followed by a close bracket of another type.

Figure 4 (a) shows the pushdown automata for the tags-only case that allows a well-formed string to be verified. Note that the only scenarios not handled in Figure 4(a) are precisely the above three scenarios: (1) a close bracket is read when the stack is empty; (2) a terminated string when the stack is non-empty; and (3) a close bracket in the string does not match the open bracket on the stack.

We have already discussed how to handle the first two scenarios. It is the third scenario that is most challenging. Note that each possible insertion operation has a symmetrically equivalent deletion operation, so for the sake of reducing the enumerated repairs we shall use a canonical form involving only deletions. Given a mismatch of types between an adjacent open-close pair, there are only five edit operations that need to be considered:

1. Delete the open bracket on the left;
2. Delete the close bracket on the right;
3. Substitute the left or right bracket to make a matching pair;
4. Substitute the open bracket on the left to a close bracket;
5. Substitute the close bracket on the right to an open bracket.

For the third alternative, we shall canonically replace the right close bracket to match the left open bracket. For the last two alternatives, the way the replacement bracket is chosen is as follows. For an open bracket substituted to a close bracket, we assign the bracket matching the next open bracket on the stack. (We only consider this alternative when the stack remains non-empty after deleting the open bracket.) For a close bracket substituted to an open bracket, we wait to assign the bracket until the first close bracket is encountered that gets paired with it (until then it is a “ghost” open) and then assign it so that the pair matches; if the string terminates before such a pairing occurs then it gets resolved to a deletion rather than a substitution.

The following claim establishes the correctness of our algorithm.

**CLAIM 3.** *By considering edit operations at only one of the following scenarios, a sequence of choices exists that leads to the optimal edit distance: (1) an empty stack when a close bracket occurs; (2) a non-empty stack when the string has terminated; and (3) an open bracket of one type adjacent to a close bracket of a different type. Furthermore, exhaustive branching to the five edit alternatives above leads to an optimal repair.*

*Proof Sketch.* Since by Claim 2, removal of well-formed substring does not affect the edit distance, one can greedily match well-formed substrings whenever possible. Thus after each edit operation, newly created well-formed substrings can be matched by an optimal algorithm. If well-formed substring removal does not return an empty string then the edit distance is non-zero, and either of the three scenarios listed in the claim must happen. Therefore, taking actions in these three scenarios is sufficient for guaranteeing optimality. Also, our exhaustive branch and bound algorithm considers all possible repairs at these three scenarios, ensuring that there exists a branch (edit script) that corresponds to the repairs made by an optimal algorithm.  $\square$

We consider two heuristics for choosing from these alternatives, the first of which makes a greedy decision based on local information and the second of which is based on non-local information but ignores interleaving between bracket types:

- **MaxBenefit:** At each mismatch, consider all five alternatives and take the one that enables the largest well-formed substring to be pruned (the size of which is the benefit). The time to test these alternatives is amortized: an alternative resulting in a larger number of matched brackets takes longer time but also advances that much further along, requiring in total linear space and time. When  $|T| = 1$ , MaxBenefit finds an optimal cost repair.
- **MinCost:** Precompute the imbalance for each bracket type subsequence (similar to the  $|T| = 1$  case) as follows. Let  $\sigma_a(s)$ , for  $a \in T$ , denote the subsequence of  $s$  containing brackets  $a$  or  $\bar{a}$ . For each  $a \in T$  and each suffix of  $\sigma_a(s)$ , we find the remaining subsequence after matching pairs elimination. Suppose the result for  $\sigma_a(s)$  is  $\bar{a}^i a^j$  and that there are currently  $k$  open brackets  $a$  on the stack. Then the number of unbalanced brackets in  $\sigma_a(s)$  is  $|k - i| + j$ . Taking all the subsequences  $\sigma_a(s)$ , for each  $a \in T$ , the minimum number of edit operations to well-formedness (via substitutions) is  $\lceil (\sum_{a \in T} |k_a - i_a|)/2 \rceil + \lceil (\sum_{a \in T} j_a)/2 \rceil$ . This gives a lower bound on the edit distance. So at each mismatch, the alternatives are considered in turn and the one which best improves the lower bound is chosen. This strategy can be done in linear space and time since the imbalance counts (for each subsequence suffix) can be precomputed and stored globally. When  $|T| = 1$ , MinCost also finds an optimal cost repair.

Unfortunately, both of these heuristics can result in approximations of the edit distance with a performance ratio that is linear in  $n$ . For example, the string  $aaaaaabbb\bar{a}\bar{a}\bar{a}\bar{a}\bar{a}$  could result in 8 edit operations with MaxBenefit if the wrong alternative among ties is chosen at each mismatch (at mismatch of type  $\bar{b}\bar{a}$ , substitute  $\bar{a}$  to  $a$ , instead of substituting  $b$  to  $\bar{b}$ ), and the string  $abcde\bar{a}\bar{e}\bar{d}\bar{c}\bar{b}$  could result in 8 edit operations with MinCost if the wrong alternative among ties is chosen at each mismatch. To mitigate this, all of the ties can be maintained in a queue and tried as part of a branch-and-bound algorithm; we call these variants MaxBenefit ++ and MinCost ++, respectively.

Interestingly, MinCost performs well on the hard input for MaxBenefit and MaxBenefit performs well on the hard input for MinCost: the first string can be repaired in 2 operations using MaxBenefit and the second in 2 operations using MinCost. Therefore, we consider hybrid strategies which combine MaxBenefit and MinCost in various ways to complement each other. In particular, we consider the following three hybrids.

- **Conservative:** Try all the choices in the union that MaxBenefit or Min-Cost gives.

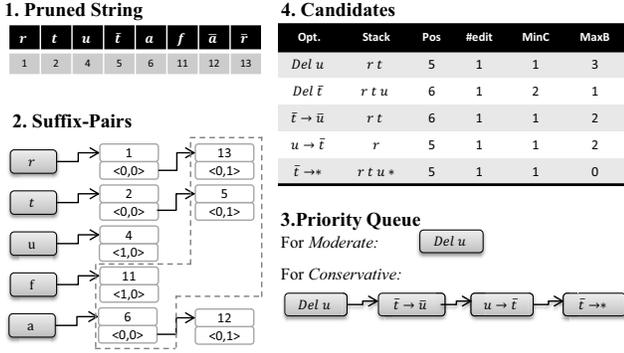


Figure 3: Illustration of Branch-and-Bound Methods

- **Moderate:** Try all the choices in the multiset intersection of the choices that Max-Benefit or Min-Cost gives; if the intersection is empty, then try all the choices from their union.
- **Liberal:** Select one choice at random from the multiset intersection of the choices that Max-Benefit or Min-Cost gives; if the intersection is empty, then select one choice at random from their union.

The above three strategies trade off between accuracy, the number of repairs returned and running time. Conservative has the highest accuracy and returns the maximum number of alternatives, whereas Liberal runs the fastest. Moderate seems to have the best trade off between accuracy and running time (see Section 5).

We note that none of these three approaches return all possible repairs within a given edit distance. The number of such repairs can be exponential in the number of errors and presenting the user with all repairs based on the many ways to resolve the errors can be overwhelming. Instead, we try to consolidate the repairs with trivial differences into representative ones: (1) while substitutions can be repaired by either changing the left tag to match the right or vice versa, we canonically choose the former; (2) while an unmatched tag can be resolved by either deleting that tag or inserting its counterpart, we canonically choose the deletion option; and (3) while a run of unmatched open (respectively, close) tags can be corrected with many different embeddings, we canonically report the maximum or minimum depth embedding. Consolidating multiple repairs by such representatives helps to provide more variety in a small set of repairs returned to the user.

### 3.3 Implementation for Branching Strategies

Figure 3 illustrates how the branch-and-bound algorithm works on the tokenized string  $rtu\bar{t}an\bar{n}f\bar{a}\bar{r}$  from Figure 1(a). Two global structures, bracket list and suffix pairs, are preprocessed in a single scan to assist the procedure. The bracket list contains the tokenized brackets and their index positions after well-formed substring pruning. The suffix pairs are constructed based on the suffix of the string starting from position  $pos$ , and for the subsequence  $\sigma(x)$  corresponding each type of bracket  $x$ , where a pair  $\langle i, j \rangle$  represents the remaining  $\bar{x}^i x^j$  after well-formed substring pruning. For example, bracket  $r$  has two suffix pairs:  $\langle 0, 0 \rangle$  at position 1 and  $\langle 0, 1 \rangle$  at 13.

The first mismatch occurs at  $\bar{t}$  (position 5). For each of the five edit alternatives, we list the stack state, the string position after the edit is applied and the total number of edits incurred for the repaired prefix. In addition, we show the MaxBenefit and MinCost values for each alternative. Taking option  $Del\ u$  for instance,  $t$  and  $\bar{t}$  also get matched, so the total benefit value is 3. To get the

MinCost value, we consider the suffix pairs surrounded by dashed lines. With  $r$  and  $t$  in the stack,  $Del\ u$  reduces the imbalance from the suffix pair  $\langle 1, 0 \rangle$  at position 4 for  $u$  by 1; the pairs  $\langle 0, 1 \rangle$  at 13 for  $r$  and  $\langle 0, 1 \rangle$  at 5 for  $t$  get canceled out by the stack; and of the remaining pairs,  $\langle 0, 0 \rangle$  at 6 for  $a$  gives 0 (since brackets can potentially match) and  $\langle 1, 0 \rangle$  at 11 for  $f$  gives 1, resulting in a total value of 1.

We show the resulting priority queues for Moderate and Conservative at the bottom right of Figure 3. The former contains the intersection of alternatives where MinCost of 1 is the lowest and MaxBenefit of 3 is the highest; the latter contains the union of alternatives having MinCost of 1 or MaxBenefit of 3. Candidates are inserted into the priority queue sorted on ascending order of  $(\#edit + MinCost)$ ; this provides a lower bound on the eventual repair cost. By visiting nodes in this order, it is guaranteed that any fully repaired string must have edit distance no larger than the existing partial repairs will have after completion.

Following the Moderate priority queue, the next mismatch occurs at  $\bar{a}$ , for which the alternative chosen is  $Del\ f$ , with MinCost value of 0 and MaxBenefit value of 5. The final output is  $rt\bar{t}an\bar{n}\bar{a}\bar{r}$  with cost of 2. Conservative returns the same repaired string with the same edit distance but uses more space and time for the extra candidates generated.

To speed up the algorithm, we can avoid visiting candidates with the same stack state and string position but having larger cost. This can be done by hashing the candidate's stack and the (index of the) remaining string suffix; when multiple repairs are needed, the hashing function is based on the repaired prefix rather than the stack.

## 4. TAGS WITH TEXT CASE

A commonly occurring pattern for semi-structured documents, especially those used for data interchange, is that text content must occur between all matching leaf-level tags and only between such tags, that is, it cannot occur immediately before or after any intermediate (non-leaf) tag. This is the norm for JSON files as well as XML files encoding JSON. Moreover, a study of DTDs on the Web revealed that only 1% of XML data exchange documents allowed so-called *mixed content elements* (allowing both text and tags) [7]. Indeed, we can exploit this pattern to compute a more judicious repair.

Let  $\Sigma$  be some alphabet and  $W = \{w \mid w \in \Sigma^+\}$  denote a set of words that can be embedded in a semi-structured document. We assume that the input data has been preprocessed into a sequence of brackets and words (which assumes the existence of markers that tell the lexical analyzer how to distinguish between brackets and words).

*Definition 7.* A well-formed string over some bracket alphabet  $T$  with embedded text from  $W$  obeys the context-free grammar  $G_{T,W}$  with productions  $S' \rightarrow S \mid \varepsilon$  and  $S \rightarrow SS \mid xS\bar{x} \mid xw^+\bar{x}$ , for all  $x \in T$  and where all  $w \in W$ .

*Definition 8.* A well-formed bracketed language with text  $L(G_{T,W})$  over some set of words  $W$  and bracket alphabet  $T$  is the set of strings from  $(W \cup T \cup \bar{T})^*$  accepted by the grammar  $G_{T,W}$  defined above.

**Example:** Let  $W = \{w\}$  and  $T = \{a, b\}$ . Then  $abu\bar{b}\bar{a}$  is a well-bracketed string, since it can be parsed as  $S \rightarrow aS\bar{a} \rightarrow abu\bar{b}\bar{a}$  but  $ab\bar{b}w\bar{a}$  is not.  $\square$

*Definition 9.* The Bracketed Language with Text Edit Distance Problem, given a string  $s \in (W \cup T \cup \bar{T})^*$ , is to find

$\arg \min_{s'} E(s, s')$  such that  $s' \in L(G_{T,W})$ . Here we allow insertion, deletion and substitution operations on brackets but do not allow any operations on words, that is, the words in a string must remain as they are.

**Example:** The string  $abb\bar{w}\bar{a}$  has edit distance 2 and can be repaired to  $abw\bar{b}\bar{a}$  using one deletion and one insertion.  $\square$

## 4.1 Optimal Solution

As shown in Figure 2, if we consider only the bracket subsequence of the string and apply the tags-only methods from Section 3 to find a solution, the resulting repairs may not obey  $G_{T,W}$ . In fact,  $G_{T,W}$  is strictly more constrained than  $G_T$  (hence, the edit distance for tags-only lower-bounds that for tags with text). Therefore, we need to design a new algorithm.

Algorithm 2 presents the pseudocode for the dynamic program. As with the tags-only case, the algorithm runs in  $O(n^3)$  time and requires  $O(n^2)$  space. Given some substring  $s_i \dots s_j$ , the algorithm first checks if it contains a word (that is,  $s_k \in W$  for some  $k \in [i, j]$ ) and, if not, deletes  $s_i \dots s_j$ , resulting in cost  $C[i, j] \leftarrow j - i + 1$ . Otherwise, either  $s_i$  and  $s_j$  could be edited to matching brackets surrounding a well-formed substring (the  $S \rightarrow xS\bar{x}$  production); or  $s_i$  and  $s_j$  could be edited to matching brackets surrounding a sequence of one or more words, after deleting all brackets in the substring  $s_{i+1} \dots s_{j-1}$ , denoted by  $D[i+1, j-1]$  (the  $S \rightarrow xw^+\bar{x}$  production); or else the string could be broken into two adjacent well-formed substrings (the  $S \rightarrow SS$  production).

Let  $B[i, j]$  be the smallest cost of editing  $s_i$  and  $s_j$  to match. When  $j = i$ , if  $s_i \in W$ ,  $C[i, j] = C[i, i] = 2$ , else  $C[i, j] = C[i, i] = 1$ . For all other cases, while the substring  $s_i s_{i+1} \dots s_j$  contains a word, the recurrence is

$$C[i, j] \leftarrow \min \begin{cases} B[i, j] + C[i+1, j-1], \\ B[i, j] + D[i+1, j-1], \\ \min_{i \leq k \leq j-1} C[i, k] + C[k+1, j] \end{cases}$$

If  $s_i \in W$  and  $s_j \in \bar{T}$ , then  $B[i, j] = 1$ . Similarly, if  $s_j \in W$  and  $s_i \in T$ , then  $B[i, j] = 1$ . Else if,  $s_i \in W$  and  $s_j \in T$ , then  $B[i, j] = 2$ . Also, if  $s_j \in W$  and  $s_i \in \bar{T}$ , or both  $s_i$  and  $s_j$  are words, then  $B[i, j] = 2$ .

---

### Algorithm 2 Dynamic Program for Tags with Text

---

**Require:** tokenized string  $s = s_1 \dots s_n$   
**Ensure:** edit distance  $\min_{s'} E(s, s')$  where  $s'$  is well-formed

- 1: **for all**  $\ell$  from 1 to  $n - 1$  **do**
- 2:   **for all**  $i$  from 1 to  $n - \ell$  **do**
- 3:      $j \leftarrow i + \ell$
- 4:     **if**  $s_k \notin W, \forall k \in [i, j]$  **then**
- 5:        $C[i, j] \leftarrow j - i + 1$
- 6:     **else**
- 7:        $C[i, j] \leftarrow B[i, j] + D[i+1, j-1]$
- 8:        $C[i, j] \leftarrow \min(C[i, j], B[i, j] + C[i+1, j-1])$
- 9:       **for all**  $k$  from  $i$  to  $j - 1$  **do**
- 10:          $C[i, j] \leftarrow \min(C[i, j], C[i, k] + C[k+1, j])$
- 11:       **end for**
- 12:     **end if**
- 13:   **end for**
- 14: **end for**
- 15: **return**  $C[1, n]$

---

CLAIM 4. Algorithm 2 correctly finds the edit distance, given a string  $s$ , such that it is accepted by  $G_{T,W}$ .

The proof of this claim can be found in the full version.

Similar to the tags-only case, we can preprocess the string to enable faster computation by removing well-formed substrings using a stack. We must also mark the presence of each removed substring to allow local edit operations, such as surrounding matching brackets that would not be allowed otherwise. This pruning can speed up the dynamic program by shortening the input string. Unfortunately, it does not guarantee a repair with optimal edit distance. For example, the string  $aaaw\bar{a}\bar{a}aw\bar{a}\bar{a}\bar{a}$  has edit distance 1 (by replacing the second  $\bar{a}$  to  $a$ ) but the pruned string  $aw\bar{a}\bar{a}$  has edit distance 2. We show that the above dynamic program stays within at least a factor of 2 after well-formed substring pruning.

CLAIM 5. Removing well-formed substrings obtains a 2-approximation on the edit distance.

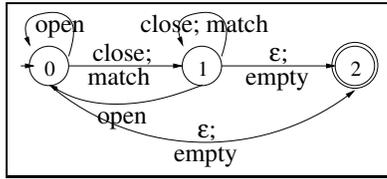
*Proof Sketch.* Define by a block a substring of brackets that appears between two consecutive occurrence of texts, or before the first text, or after the last text. Consider a new edit distance function, which is exactly similar to the original edit distance problem, except that two consecutive open brackets (close brackets) in a block can both be deleted at a cost of 1. If this new edit distance function has minimum distance  $d$ , then clearly, the optimal edit distance for the original problem cannot be more than  $2d$  (we pay 2 to delete two open (close) brackets in a block). The main insight is to consider this new function, and show that well-formed substring pruning preserves edit distance for it, and hence guarantees 2-approximation to the original problem. The proof works by doing induction on the number of errors and showing by case analysis that if an optimal algorithm does edit to a well-formed substring, then there exists an alternate edit script where the positions of the edits are outside of the well-formed substrings similar to Claim 2.  $\square$

## 4.2 Incremental Approach

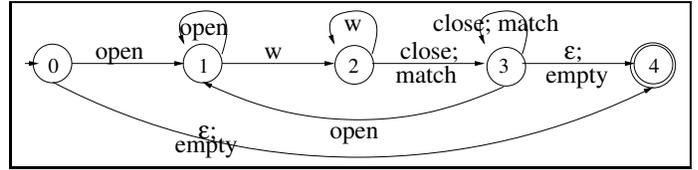
Once again, the dynamic program is intended for computing the edit distance, along with perhaps a single repair, but does not provide a natural way of enumerating multiple repairs. Here we identify the different scenarios in which repairs should be handled as well as a set of possible repairs for each of these scenarios, similar to our approach for the tags-only case. We start by giving push-down automata in Figure 4(b) that allows a well-formed string to be verified, for the tags with text case. We have shown in Section 3 that only considering the scenarios not handled by the push-down automata for tags-only case, rather than eagerly repairing the string, leads to the optimal edit distance while making the running time proportional to the number of errors rather than the size of the string. We apply the same idea to tags with text: only scenarios not covered by the automaton in Figure 4(b) are considered. We show how to deal with these scenarios in Table 1.

By contrast, there are nine scenarios based on the automaton for tags with text in Figure 4(b). Furthermore, the alternatives for open-close bracket mismatch depend on which state the mismatch occurs in. If this occurs in State 3 then we consider the following alternatives (similar to tags-only): make the open and close brackets match via substitution; pop the open bracket from the stack; delete the close bracket in the string; substitute the open bracket in the stack to close; or substitute the close bracket in the string to ghost-open (equivalently, substitute the right to match the left). However, if this occurs in State 2, then there are different options since the word(s) must be surrounded by a pair of brackets: make the open and close brackets match via substitution; insert a matching close bracket, if possible; or delete the close bracket if the next token in the string is a word (otherwise insert a matching open bracket).

Any input string can be partitioned into *blocks* of brackets separated by text. There are five additional scenarios in addition to the



(a) tags-only



(b) tags with text

Figure 4: Automata for tags-only and tags with text

State	Token	Action	Next State
0	open	push to stack and advance string	1
0	close	sub close in string to ghost open	0
0	word	delete close from string	0
0	null	insert ghost open before word	0
0			4
1	open	push to stack and advance string	1
1	close	sub open from stack to close and pop matching pair from stack, if possible	1
1		pop open from stack, if possible	1
1		sub close from string to ghost open	1
1		delete close from string	1
1	word		2
1	null	clean-up	4
2	open	insert matching close to string before open	2
2		sub open in string to matching close	2
2		delete open from string	2
2	match	pop from stack and advance string	3
2	close	sub string to match stack	3
2		insert matching close in string, if possible	2
2		if next token is word, delete string close	2
2		else push matching open to stack	2
2	word		2
2	null	insert matching close to string	2
3	open	push to stack and advance string	1
3	match	pop from stack and advance string	3
3	close	sub close in string to match stack	3
3		delete close in string	3
3		pop open from stack	3
3		sub open in stack to close and pop resulting match from stack, if possible	3
3		sub close in string to ghost open	3
3	word	insert ghost open in string (before word)	3
3	null	clean-up if non-empty stack	4

Table 1: State Transition Table for Tags with Text Case

three for tags-only: (4) a close bracket occurs immediately after an open bracket; (5) an open bracket occurs immediately after a word; (6) a word occurs immediately after a close bracket; (7) a word occurs as the first token; and (8) the string terminates in a word. For these additional scenarios, there are various edit alternatives, which are listed in Table 1.

The so-called “clean-up” phase referred to in States 1 and 3 of Table 1 is invoked if the stack is non-empty when the string terminates. In this case, the goal is to take the existing stack, paying attention to the blocks that each stack open bracket is part of, and perform the minimum number of substitutions and deletions to obtain a well-formed string. For example, suppose there are three blocks on the stack, the first with  $ab$ , the second with  $cde$  and the third with  $fg$ . By deleting  $d$  and replacing  $c$  with  $\bar{b}$ ,  $f$  with  $\bar{e}$  and  $g$  with  $\bar{a}$ , the resulting brackets are well-formed. We omit details for lack of space.

**Example:** Let  $s = a\bar{a}\bar{w}a\bar{a}w\bar{a}\bar{a}$ . Repair scenario (4) occurs after the pair  $a\bar{a}$  in the first block (since there is no text separating them), causing  $\bar{a}$  to either be replaced with an open bracket or deleted (the other two alternatives from State 1 are not possible). Scenario (4) occurs again at the next  $\bar{a}$  in the string with the same edit alternatives. Suppose we choose the substitution al-

ternative both times. Then Scenario (6) occurs after the first  $w$ , which we can repair by either inserting a close, substituting the open to a close or deleting the open. Suppose we choose to insert a close. The remaining elements will be read without problem until the string terminates, at which point the stack will be non-empty with two open brackets. At that point, they must both be deleted since they both occurred in the same block. The final repair, then, is  $aaaw\bar{a}aaaw\bar{a}\bar{a}$  with a cost of 5. Had we instead chosen the alternative to delete the close brackets in Scenario (4), the string would have been repaired to  $aw\bar{a}aaaw\bar{a}\bar{a}$  at a cost of 3, which is optimal.  $\square$

The following claim establishes the correctness of the automaton. Similar to Claim 3, the main idea is to show the states considered in our automaton are the only scenarios where repair has to be made if well-formed substrings are removed greedily. If we consider an optimal algorithm for the new edit function defined in Claim 5, then it can be shown that at each error state, repair choices considered are exhaustive and hence there exists a branch leading to optimal cost for the new function. Since, any optimal algorithm for the new edit function returns a solution within twice the minimum edit distance of our problem, the claim is established.

**CLAIM 6.** *The automaton given in Table 1 with branches to all the above edit alternatives, obtains a 2-approximation on edit distance.*

As before, for the MaxBenefit strategy, all alternatives are considered in turn and then the one resulting in the largest number of brackets that can be paired to matches is chosen. For MinCost, we employ the same cost estimation formula as the tags-only case, since it provides a lower bound; the alternatives are sorted with respect to these costs.

## 5. EXPERIMENTAL EVALUATION

This section gives a thorough experimental evaluation of methods for both Tags-Only and Tags With Text problems, against the *Dynamic Programming* (DP) and *branch-and-bound* algorithms.

### 5.1 Set-up

All algorithms were implemented in Java and executed on a server with a Quad-Core AMD Opteron(tm) Processor 8356@1 GHz and 128 GB RAM running Centos 5.8. We used the following three real data sets:

- *BGP*<sup>4</sup> real-time routing information provided by BGPmon. We used a portion of the stream output over some time interval.
- *Tree Bank*<sup>5</sup> annotated linguistic text, with average depth 7.8 and max depth 36. We extracted random subtrees with max-depth no less than 20 and merged them together.

<sup>4</sup><http://bgpmon.netsec.colostate.edu/>

<sup>5</sup><http://www.cis.upenn.edu/~treebank/>

- *Web Repository*<sup>6</sup> a Web repository of XML documents containing 180,000 files, 10% of which are reported with begin and end tag mismatched.

Both *Tree Bank* and *BGP* normally satisfy the grammar  $G_{T,W}$ . For the XML Collections, we picked two sets of 80 documents which also satisfy the grammar  $G_{T,W}$ . For the Tags-only methods, only tag subsequences are retained and for tags with text the subsequence of tags and text was retained; all attributes were removed.

**Error Model:** We chose from among six different operations, with equal probability, to inject errors into a given well-formed string; the error types are listed in Table 2. The position  $i$  of each operation was either chosen uniformly (without replacement) in  $[1, n]$  or non-uniformly following a normal distribution about a randomly-chosen “center” with some standard deviation but keeping within  $[1, n]$ .

Table 2: Types of Injected Errors

Operation	Description
Delete( $i$ )	delete the tag at the $i$ -th position
Insert( $i, a$ )	insert tag $a$ to the $i$ -th position, where $a$ is randomly chosen from $(T \cup \bar{T})$
Swap( $i, j$ )	swap the tag located at the $i$ -th position with the one at the $j$ -th
Flip( $i$ )	change the $i$ -th tag to close (resp., open) if it is open (resp., close)
Sub( $i, a$ )	substitute the $i$ -th tag with $a$ , where $a$ is randomly selected from $(T \cup \bar{T})$
DeepInsert( $a, h$ )	insert tag $a$ into some position $i$ having $depth(i) > h$ , where $a$ is randomly selected from $(T \cup \bar{T})$

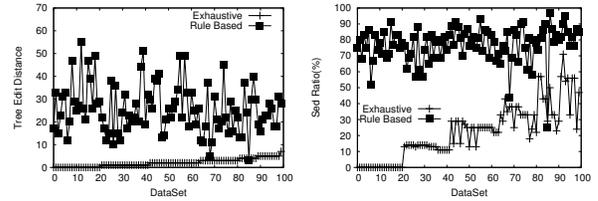
## 5.2 Effectiveness of Edit Distance Approach

To evaluate the goodness of our edit distance based approach, we started from a well-formed string  $s$ , injected errors randomly to obtain string  $s'$ , ran our methods to find a repair  $s''$ , and compared  $s''$  against the original well-formed string  $s$ . Since the original string and the repair are both well-formed, we are now able to use approximate tree matching algorithms to evaluate the goodness of our repairs. For this, we used the recently developed RTED algorithm from [21], which uses the standard tree edit operations [31]: delete a node and connect its children to its parent, maintaining the order; insert a new node between an existing node,  $v$ , and a consecutive subsequence of  $v$ 's children; and rename the label of a node. We compared our approach against the rule-based heuristics mentioned in Section 1 for dealing with open-close tag mismatches: 1) delete the open tag; 2) delete the close tag; and 3) substitute the close tag to match. In the following experiments, we present the best of these three heuristics.

We used a data set of 100 strings from *Tree Bank* consisting of 1,000 tags and around 8 errors, with uniformly distributed errors. (The results with normally distributed error are similar, so we omit them for lack of space.) Figure 5(a) shows the tree edit distance using both the two methods for the Tags-with-Text case, where the x-axis indexes all of the 100 data sets sorted by the tree edit distance of Exhaustive repair, which is the branch-and-bound algorithm that tries all five choices at each branch point. There were 20 strings for which our method obtained  $d_{TED}(s, s'') = 0$ , a complete reversion of the string compared to the original, and over 60% of the strings had  $d_{TED}(s, s'') \leq 3$ . In contrast, the rule-based heuristics obtained strings with average tree edit distance 25.

In addition, we also did the following comparison. Starting from the original (well-formed) input string  $s$ , we injected errors to obtain string  $s'$  and computed the *string* edit distance between  $s$  and

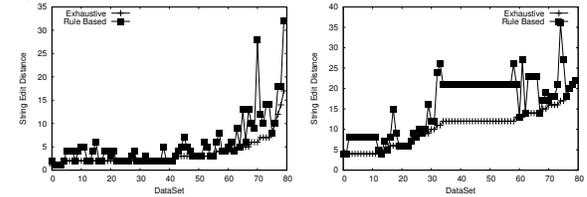
<sup>6</sup><http://data.politicalmashup.nl/xmlweb/>



(a) Tree Edit Distance

(b) String Edit Ratio

Figure 5: Goodness of Exhaustive and Rule-Based repairs (*Tags With-Text*)



(a) Tags-only

(b) Tags with Text

Figure 6: String Edit Distance with Real Errors

$s'$ . Then we repaired the string to obtain string  $s''$  and computed the string edit distance between  $s'$  and  $s''$ . Finally we calculated  $\frac{d_{SED}(s, s'')}{d_{SED}(s, s') + d_{SED}(s', s'')}$  as a measure of how well the error was “undone”, where 0 means exact reconstruction and 1 means the repair resulted in no improvement. We compared our approach with the rule-based heuristics using this measure on real data. Figure 5(b) plots this ratio for both methods. The rule-based heuristic ranges from 50% to 100%, with an average of 80% while Exhaustive had an average ratio of 20% with about 20 strings having ratio 0. In fact, Exhaustive beat rule-based on all the 100 strings except one. From these experiments, we conclude that our edit distance based approach succeeds in reverting the strings towards their original form, far more than the rule-based heuristics.

We also tested our approach using data from *Web Repository*, which contains real malformed XML data. Since no gold standard of what was intended by the creator of any of these XML documents was available to evaluate goodness of repair, our experiments compare the number of operations to obtain a well-formed string using our methods with that obtained using the rule-based heuristics.

Figure 6 presents the results for both Tags-only and Tags-with-Text cases on 80 malformed XML documents chosen randomly with number of tags ranging from tens to thousands. Figure 6 presents the results. With few errors in the string, our approach was only marginally better than the rule-based. However, there was up to an order of magnitude difference for strings requiring at least 10 edit operations to repair. This trend was evident for both cases shown in Figure 6(a) and Figure 6(b), respectively.

In the following sections, we study the performance of our algorithms by measuring *Average Running Time* and *Average Edit Distance*.

## 5.3 Tags-Only, Single Repairs

We compared methods that are designed to return a single repair, including DP as well as MaxBenefit, MinCost and Hybrid, where Hybrid randomly picks one choice either given by MinCost or by MaxBenefit. The rule-based heuristics were used as a baseline here. Finally, we tried five trials of Random, which randomly chooses one from the five alternatives at each tag mismatch, and reports the lowest edit distance among these. By default, well-formed substring pruning is applied to speed up the methods.

Figure 7 shows the edit distance and running time as a function

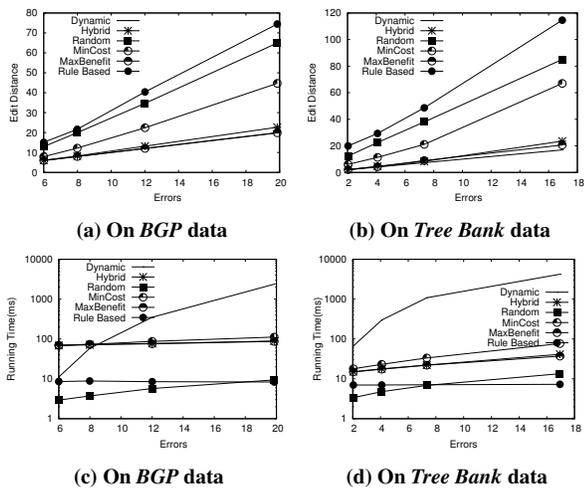


Figure 7: Single Repair, Error Number (TagsOnly)

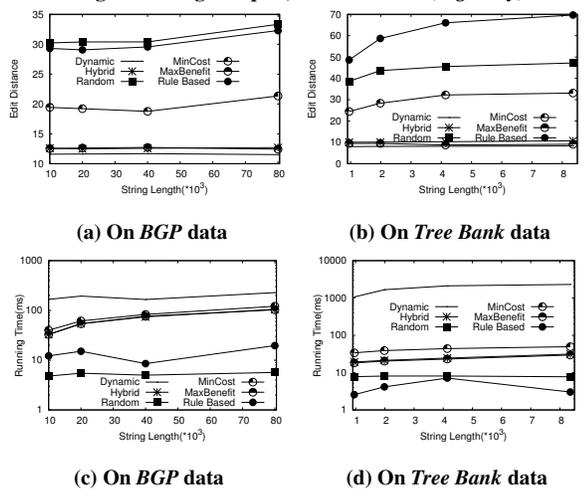


Figure 8: Single-Repair, String Length (TagsOnly)

of the number of errors. For *BGP*, the errors ranged from 6 to 20 with initial string length fixed at 40,000; after well-formed substring pruning, the string length was significantly reduced to the range [20, 160]. For *Tree Bank*, the error ranged from 2 to 16 with string length fixed at 4,000; after pruning the length was reduced to [40, 280].

While DP gave the smallest edit distance (it is optimal), it was also the slowest in almost all cases. The running time of DP increased significantly with the number of errors, even with well-formed substring pruning: the string length after pruning increased from 20 to 160 on *BGP* and from 40 to 280 on *Tree Bank*. Rule-Based, in contrast, ran the fastest but gave the highest edit distance. Interestingly, Rule-Based was even less accurate than Random (recall that Random chooses from among five alternatives while Rule-Based makes a single deterministic choice), which in turn also affected its running time performance due to the additional edits. The edit cost of MaxBenefit was close to optimal, its inaccuracy growing with increasing errors but at a very slow rate; at the same time, its running time was 1-2 orders of magnitude faster than that of DP. MinCost, on the other hand, was no faster but much less accurate than MaxBenefit. Hybrid, which integrates both heuristics, is slightly less accurate than MaxBenefit but slightly faster.

Figure 8 demonstrates the scalability with respect to string length, with the number of errors fixed at around 12 on *BGP* and 8 on *Tree Bank*. Due to well-formed substring pruning, the average

string length was reduced to around 60 on *BGP* and 150 on *Tree Bank* from all the initial string lengths. Hybrid and MaxBenefit follow DP closely in accuracy and outperform the latter in running time, as much as 1-2 orders of magnitude. As proved in Sect. 3.1,

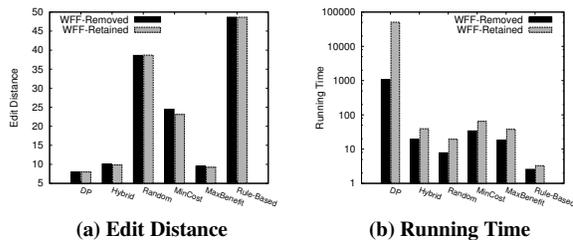


Figure 9: Well-formed Substring (TagsOnly)

well-formed substring pruning, does not affect the edit distance of DP. However, it does affect running time. We ran a set of experiments on *Tree Bank* with average string length of 1,000 and 8 errors on average. The bars in Figure 9 show the difference in running time and edit distance with and without such pre-processing. The edit distance for DP and Rule-Based stays the same, but there is a slight difference in other methods, which can be explained by randomization. For DP the running difference is quite noticeable (from 1066ms to 50510ms) while, for others, the increase in running time is small; such increase is largely brought by the cost in building the global suffix stack.

## 5.4 Tags-Only, Multiple Repairs

We compared various branch-and-bound methods: Exhaustive, Conservative, Moderate, Liberal, MinCost++ and MaxBenefit++. The key difference between these methods is the number of alternatives tried at each branch, where there is an inherent trade-off between accuracy and running time.

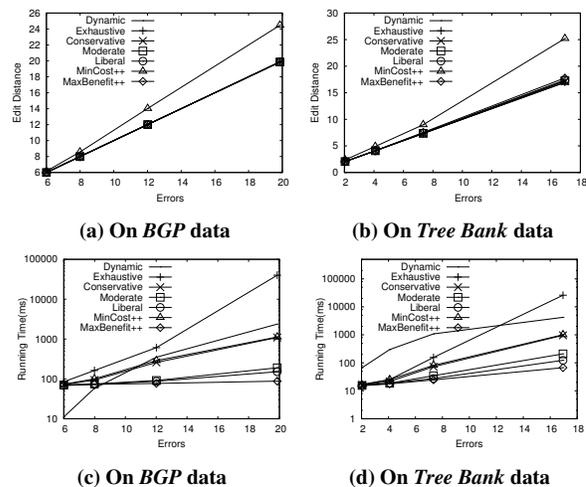
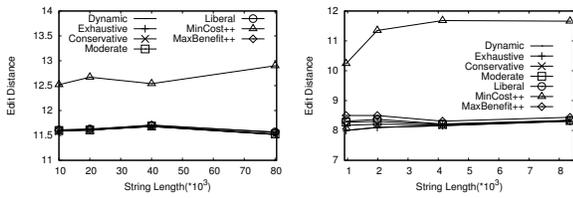


Figure 10: Multi-Repairs, Error Number (TagsOnly)

**Single-Repair Performance** We begin by showing the results for single repair, i.e.,  $K = 1$ , using DP as a baseline.

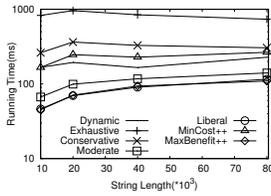
Figure 10 shows performance versus error number ranging from 6 to 20 on *BGP* and 2 to 16 on *Tree Bank*. As expected, Exhaustive gives optimal edit distance while MinCost++ is the least accurate one. In Figure 10(a) to (d), all methods except MinCost++ and Exhaustive are almost as accurate as DP but much faster. MaxBenefit++ beats all other methods in running time, but is less accurate.

Figure 11 shows performance versus string length ranging from 10,000 to 80,000 (with roughly 12 errors) on *BGP*, and 1,000 to

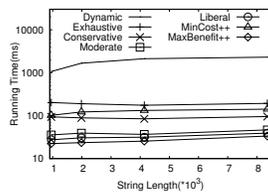


(a) On BGP data

(b) On Tree Bank data



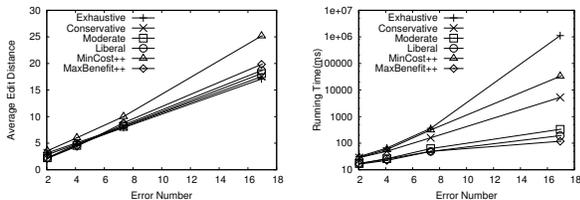
(c) On BGP data



(d) On Tree Bank data

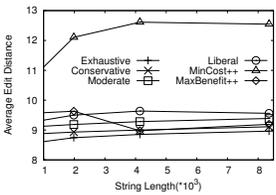
Figure 11: Multi-Repair, String Length (TagsOnly)

8,000 (with roughly 8 errors) on *Tree Bank*. After pre-processing, the string sizes were greatly reduced to 60 on *BGP* and 130 on *Tree Bank*. When a string is short with many errors, DP wins; otherwise, Exhaustive is faster. In general, the branch-and-bound methods were not greatly affected by string length and perform well when the number of errors is small.

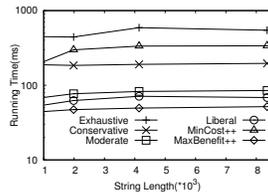


(a) Edit distance vs #errors

(b) Running time vs #errors



(c) Edit Distance vs length



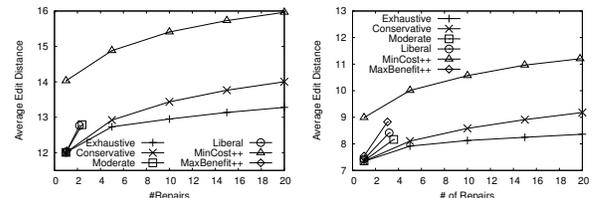
(d) Running time vs length

Figure 12: top-5 Repairs (TagsOnly)

**Multi-Repairs Performance** Figure 12(a) and (b) illustrates the performance for finding 5 repairs when the number of errors increases from 2 to 16 on *Tree Bank* (with string length fixed at 4,000); the results on *BGP* showed similar trends and are omitted for lack of space. Figure 12(c) and (d) illustrates the performance when string length grows from 1,000 to 8,000 on *Tree Bank*. With the well-formed substring removed, the string length decreases significantly to 130 on *Tree Bank*.

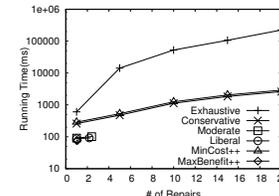
Exhaustive, Conservative and MinCost++ are 10 times slower than Moderate, Liberal and MaxBenefit++ since the former methods return more repairs than the latter do. MinCost++ is less accurate than Exhaustive and Conservative, but is faster on *BGP*; on *Tree Bank* MinCost++ is less accurate and comparable to Exhaustive in running time.

**K-Repairs Performance** We evaluated the performance of finding up to  $K$  repairs, for  $K \in \{1, 5, 10, 15, 20\}$ , with string length 40,000 for *BGP* and 4,000 for *Tree Bank*. Note that not all methods were able to obtain  $K$  repairs. Moderate, Liberal and

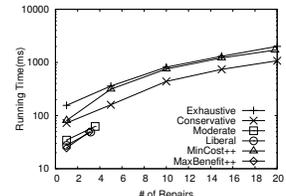


(a) On BGP data

(b) On Tree Bank data



(c) On BGP data

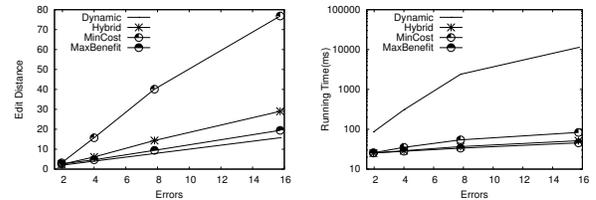


(d) On Tree Bank data

Figure 13: top-k, Scalability (TagsOnly)

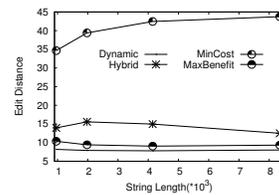
MaxBenefit++ run faster by aggressively pruning; therefore, they result in fewer total repairs. Figure 13 shows that these three methods were unable to return more than 5 repairs. Only Exhaustive, Conservative and MinCost++ obtained up to 20 repairs. With more nodes visited, Exhaustive returned repairs lower in average edit distance but requires running time. On *Tree Bank*, MinCost++ is worse in both average edit distance and running time, which means **MinCost** prunes off some nodes low in edit distance, leading to a longer edit path and larger search space. For methods where there are enough repairs, the running time grows linearly in  $K$ .

## 5.5 Tags with Text, Single Repair

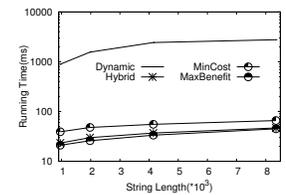


(a) Edit distance vs #errors

(b) Running time vs #errors



(c) Edit distance vs length



(d) Running time vs length

Figure 14: Single-Repair (TagsWithText)

We ran experiments using the dynamic program for repairing documents with text (satisfying  $G_{T,W}$  rather than  $G_T$ ) as well as the following analogues of tags-only methods: MaxBenefit, MinCost and Hybrid. Figure 14 presents performance as a function of number of errors, ranging from 2 to 16 (with fixed string length 4,000) on *Tree Bank*. (The trends with *BGP* are very similar and omitted for lack of space.) DP again was the slowest while MinCost is the least accurate. MaxBenefit was again both more accurate and faster than MinCost and Hybrid was slightly faster but less accurate than MaxBenefit. With the well-formed substring removed, the average string size from 4,000 to 60  $\sim$  500 on *Tree Bank*. The running time for DP grew quickly, due to the increase

in string length, while other methods are less affected by errors. MaxBenefit approximated DP well in edit distance and was faster by up to two orders of magnitude.

## 5.6 Tags with Text, Multiple Repairs

**Single-Repair Performance** We looked at the performance versus number of errors, ranging from 6 to 20 on *BGP* (string length 40,000) and from 2 to 16 on *Tree Bank* (string length 4,000); the trends were very similar to that of the Tags-only case so we omit the graphs for lack of space. The edit costs were close to optimal for all methods except MinCost++, especially on *BGP*. MaxBenefit++ was the fastest and DP was the slowest with one exception: when error number is 20 on *BGP* (where the string length is very small after pruning). On *Tree Bank*, when the error number is 16 and the string length is 500, Exhaustive still beats DP by an order of magnitude. In general, the advantages of branch-and-bound methods are seen with strings of large sizes and few errors. We also looked at the scalability versus string length with roughly 12 errors for *BGP* and 8 errors for *Tree Bank*. Not surprisingly, Conservative was the second best after Exhaustive. MaxBenefit++ was superior to MinCost++ on *Tree Bank* but not on *BGP*, which shows that the heuristics they're based on are complementary. After pruning, the string length remains around 150 for *BGP* and 300 for *Tree Bank*, which explains the stability in running time for DP on both datasets. Nonetheless, DP is slower than Exhaustive by two orders of magnitude on *Tree Bank*.

**Multi-Repair Performance** We looked at the performance and scalability of branching methods to get 5 repairs. On *Tree Bank*, Conservative beats MinCost++ on both accuracy and speed, which shows the MinCost heuristic does not work well in some cases as few of its branches led to low-cost repairs. The constancy in string length after pruning is the main reason why the running time for both datasets are fairly constant with increasing string size. The average edit distance of Moderate, Liberal and MaxBenefit++ seem smaller than even Exhaustive when string length equals to 8,000; however, this is partly due to them finding no more than 3 repairs. Since these trends were similar to that of Tags-only multi-repairs, we omit the graphs for brevity.

**K-Repair Performance** We issued queries to find  $K$  repairs for  $K$  between 1 to 20; the trends were similar to Tags-only  $K$ -repair performance so we omit the graphs for lack of space. The methods that prune more aggressively failed to return as many as  $K$  repairs in some instances. Only MinCost++, Conservative and Exhaustive were capable of returning  $K$  repairs. Exhaustive gave the smallest average edit distance but at a much higher running time; Conservative was comparable to MinCost++ but retrieved repairs with smaller average edit distance, especially on *Tree Bank*. With increasing  $K$ , the average edit distance grew slowly while the running time grew linearly, which indicates some degree of scalability.

## 6. CONCLUSIONS

Future work includes considering more complicated edit distances, including approximate string matching for tag labels as well as more general sequence alignment operations and allowing additional operations such as swaps and moves. In addition, there are other cues in a document, even when the grammar is unknown, that can be leveraged for a more judicious repair including attributes, text correlations, etc.

## 7. REFERENCES

- [1] A. V. Aho, I. G. Peterson: A Minimum Distance Error-Correcting Parser for Context-Free Languages. *SIAM J. Comput.* 1(4): 305-312 (1972)

- [2] G. Beskales, I. F. Ilyas and L. Golab: Sampling the Repairs of Functional Dependency Violations under Hard Constraints. *PVLDB* 3(1): 197-207 (2010)
- [3] G. J. Bex, F. Neven, T. Schwentick, K. Tuyls: Inference of Concise DTDs from XML Data, *VLDB 2006*: 115-126
- [4] G. J. Bex, F. Neven, S. Vansummeren: Inferring XML Schema Definitions from XML Data, *VLDB 2007*: 998-1009
- [5] U. Boobna and M. de Rougemont: Correctors for XML Data. *XSym 2004*: 97-111
- [6] P. Bohannon, M. Flaster, W. Fan and R. Rastogi: A Cost-Based Model and Effective Heuristic for Repairing Constraints by Value Modification. *SIGMOD 2005*: 143-154
- [7] B. Choi: What are real DTDs like? *WebDB 2002*: 43-48
- [8] G. Cobena, S. Abiteboul, A. Marian: Detecting Changes in XML Documents. *ICDE 2002*: 41-52
- [9] T. Dalamagas, T. Cheng, K. Winkel, T. K. Sellis: A Methodology for Clustering XML Documents by Structure. *Inf. Syst.* 31(3): 187-228 (2006)
- [10] W. Fan, J. Li, S. Ma, N. Tang and W. Yu: Interaction between record matching and data repairing. *SIGMOD 2011*: 469-480
- [11] M. N. Garofalakis, A. Gionis, R. Rastogi, S. Seshadri and K. Shim: DTD Inference from XML Documents: The XTRACT Approach. *IEEE Data Eng. Bull.* 26(3): 19-25 (2003)
- [12] E. M. Gold: Language Identification in the Limit, *Information and Control* 10(5): 447-474 (1967)
- [13] S. Guha, H. V. Jagadish, N. Koudas, D. Srivastava, T. Yu: Approximate XML Joins. *SIGMOD 2002*: 287-298
- [14] S. A. Greibach: The Hardest Context-Free Language, *SIAM J. Comput.* 2(4): 304-310 (1973)
- [15] S. Grijzenhout and M. Marx: The quality of the XML web. *CIKM 2011*: 1719-1724
- [16] L. Lee: Fast context-free grammar parsing requires fast boolean matrix multiplication, *J. ACM* 49(1): 1-15 (2002)
- [17] V. I. Levenshtein: Binary codes capable of correcting deletions, insertions, and reversals. *Soviet Physics Doklady* 10(8): 707-710 (1966).
- [18] F. Magniez, C. Mathieu and A. Nayak, Recognizing well-parenthesized expressions in the streaming model, *STOC 2010*: 261-270
- [19] G. Myers: Approximately Matching Context-Free Languages. *Inf. Process. Lett.* 54(2): 85-92 (1995)
- [20] G. Navarro. A guided tour to approximate string matching. *ACM Comput. Surv.*, 33(1):31-88 (2001)
- [21] M. Pawlik, N. Augsten: RTED: A Robust Algorithm for the Tree Edit Distance. *PVLDB* 5(4): 334-345 (2011)
- [22] H. Samimi, M. Schaefer, S. Artzi, T. Millstein, F. Tip and L. Hendren: Automated Repair of HTML Generation Errors in PHP Applications Using String Constraint Solving. *Int'l Conf. Software Engineering 2012*
- [23] L. Segoufin and C. Sirangelo: Constant-Memory Validation of Streaming XML Documents Against DTDs. *ICDT 2007*: 299-313
- [24] L. Segoufin and V. Vianu: Validating Streaming XML Documents. *PODS 2002*: 53-64
- [25] S. Staworko and J. Chomicki: Validity-Sensitive Querying of XML Databases. *EDBT Workshops 2006*: 164-177
- [26] N. Suzuki: Finding an optimum edit script between an XML document and a DTD. *SAC 2005*: 647-653
- [27] D. Sankoff, J. B. Kruskal: Time Warps, String Edits and Macromolecules *Theory and Practice of Sequence Comparison, 1983*. Addison Wesley
- [28] E. S. Tabanao, M. M. T. Rodrigo, M. C. Jadud: Predicting at-risk novice Java programmers through the analysis of online protocols. *ICER 2011*: 85-92
- [29] Y. Wang, D. J. DeWitt, J. Cai: X-Diff: An Effective Change Detection Algorithm for XML Documents. *ICDE 2003*: 519-530
- [30] R. A. Wagner and M. J. Fischer. The string-to-string correction problem. *Journal of the ACM*, 21(1):168-173 (1974)
- [31] K. Zhang and D. Shasha: Simple Fast Algorithms for the Editing distance between Trees and Related Problems. *SIAM J. Computing*, 18(6): 1245-1262 (1989)